### 2021 IBM Systems



### Jean-Armand Broyelle

ML/DL Cognitive Systems Technical Leader IBM Garage for Systems - Montpellier France

age for Systems - Montpellier France

Improve DeepLearning Analysis of Large Images with IBM WML-CE



### Deep learning NN Training is memory constrained



- GPUs have limited memory
- Neural networks are growing deeper and wider

• Amount and size of data to process is always growing to improve AI solutions accuracy

### **Dimension problem with Medical Images**

#### How much can a NVIDA Tesla V100 (32Gb) take in ?

- ResNet-101, batch size=32, it can take in images of 512\*512\*3 → 750kB \* 32 → 24 MB
  ResNet-101, batch size=1, it can take in image of 3880\*3880\*3 → 43 MB
- 3D ResNet-101, batch size=32, it can take in images of 92\*92\*42\*1  $\rightarrow$  339kB \*32  $\rightarrow$  10 MB
- 3D ResNet-101, batch size=1, it can take in image of 577\*577\*42\*
   → 13MB

### But Typical Resolution of Medical Image are LARGE

- Chest radiograph : 4000\*5000 uint16 → 305MB
- Computed tomography : 512\*512\*50 uint16  $\rightarrow$  200MB
- Low-dose lung CT: 512\*512\*500 uint16 → 2GB
- Digital Whole Slide Image : 100,000\*50,000\*3 uint8  $\rightarrow$  111GB

# Current approaches to deal with NN models size constraints with Large images



#### Patch/tiles-based methods



- $\rightarrow$  Information loss
- $\rightarrow$  Low detection of tiny detailled objects

- $\rightarrow$  Low detection of objects with different scaling
- $\rightarrow$  Complex post-processing

# WML-CE introduce Large Model Support

Train larger images, more complex models

#### Traditional Model Support

Limited memory on GPU forces trade-off in model size / data resolution

#### Large Model Support

Use system memory and GPU to support more complex and higher resolution data



Leveraging NVLink and CPU-GPU memory coherence enables larger and more complex models

Improves model accuracy with more images and higher resolution images

### Gpu memory usage



## Model training in gpu memory



# Model training with tensor swapping



System memory

### What's possible with Large model support?

- GoogleNet: 10x image resolution
- Resnet50: 10x image resolution
- Resnet152: 10x image resolution
- 3D-Unet brain MRIs: 5x image resolution
- DeepLabV3 2D : 10 x image resolution

## 1<sup>st</sup> use-case : 3D-UNet image segmentation 3D Dicom images



- 3D-Unet generally has high memory usage requirements
  - 3D Dicom image size > 100 MB
- International Multimodal Brain Tumor Segmentation Challenge (BraTS)
- Existing Keras model with TensorFlow backend
- Experiment with TensorFlow Large Model Support (TFLMS)







# Effect of 2x resolution on dice coefficients

(higher is better)



© Copyright IBM Corporation 2021

### **Overhead of Large Model Support with NVLink 2.0**

#### 4.5 4 3.5 Millions of voxels / seconds 7 1.5 7 1.5 Fits in memory 0.5 0 50 0 100 150 200 250 300 350 400 450

#### 3D U-Net on Power9 with 32GB NVIDIA V100

MRI resolution (cubed)

### 2<sup>nd</sup>) Digital Whole Slide Image (WSI) ResNet Classification Experiment

•Generated by slide scanner

•Resolution can be up to 200,000 \* 100,000 pixels ( 20 Billion )  $\rightarrow$  447 GB



# Curent approach: Two-Level AI Model for Cancer Detection on Whole Slide Image



Divide WSI into patches

Patch-level model (>10M Patches) Background, Benign, Cancer Classification accuracy : 98%

Slide-level model 260 Training, 100 Testing Classification Accuracy : 97%

Benign or Nasopharyngeal Cancer?

Ground Truth : Cancer, Normal Tissue Shadowed area : Cancer predicted by Al

### Overhead of Large Model Support with NVLink 2.0 ResNet50 and LMS on AC922 V100 32GB

#### **ResNet50 on Power9 with 32GB NVIDIA V100**



### Patches versus whole image



#### Epoch throughput

#### Training time to reach a given acccuracy



## **3td ) NLP : train Bert Large with LMS**

### Requirements and limitations

- High compute
- Training made
   on 16 Google
   Cloud TPUs
   (64 chips)
- High memory
- Requires
   100+ GB for
   training
   [TODO verify / graph ?]



#### MLPerf 0.6 Benchmark Cloud TPU v3 vs. Nvidia DGX-2h with GPUs

Impossible to train  $BERT_{Large}$  on best GPUs (32GB) without reducing network and parameters (batch size = 1), which harms the final accuracy.

IBM Systems Technical University © Copyright IBM Corporation 2021



pyright IBM Corporation 2021

### **Bert Large with LMS : Results**

Goals:

 Prove that LMS allows to train or finetune BERT-Large on a single GPU BATCH SIZE = 24 BERT-Large BERT-Base 0 20 40 60 80 100 120 GPU = RAM

MEMORY REQUIRED PER MODEL (IN GB)

 Reproduce accuracy achieved by Google AI (on Cloud TPUs)



### 4<sup>th</sup>) Automotive Driving : PSPNET training against HD images with LMS

Large images allow to see far & anticipate.







# **Tensorflow Large Model implementation overview**



### **PyTorch LMS - Usage and Tuning**

# Extends PyTorch's torch.cuda package to provide the following control and tuning interfaces:

torch.cuda.set\_enabled\_lms(enable)

Enable/disable Large Model Support. Parameters: **enable** (bool): desired LMS setting.

torch.cuda.set\_limit\_lms(limit)

Sets the allocation limit (in bytes) for LMS. Parameters: **limit** (int): soft limit on GPU memory allocated for tensors.

torch.cuda.set\_size\_lms(size)

Sets the minimum size (in bytes) for LMS. Parameters: **size** (int): any tensor smaller than this value is exempt from LMS reuse and persists in GPU memory.



### **PyTorch Large Model Support - Consideration**

- Ensure active tensors are resident in device memory during computations. Maintain set of inactive tensors eligible for eviction.
- Enhance GPU allocator to allow eviction of inactive tensors as an alternative to new device memory allocation.
- All tensors, including weights/biases, are eligible for swapping.
- No graph to analyze, but PyTorch efficiently manages tensor lifespans.
- Implementing heuristics which are informed by prior iterations is a challenge:
  - (Most) tensors do not persist across iterations
  - Conditional control flow allows iterations to differ from one to the next
- Next memory limitations to consider :
  - RAM capacity (to avoid system swapping)
  - CudNN : single tensor indexing is based on INT32 → the size of one single tensor can't exceed 2Giga elements of tensor 'datatype'



### **PyTorch LMS - Statistics**

# Extends PyTorch's torch.cuda package to provide the following statistics interfaces:

# Memory active: torch.cuda.memory\_active() torch.cuda.max\_memory\_active() torch.cuda.reset\_memory\_active()

#### # Memory swapped: torch.cuda.memory\_reclaimed() torch.cuda.reset\_memory\_reclaimed()

## Conclusion

- Dimension problem with large image for AI.
  - → Large Training images provide better models but hit GPU limit on traditionnal h/w
- > Classical ways to overcome dimension problems could limit model performances.







- > WML-CE Large Model Support is a unique differentiator that address this issue.
  - Native input image resolution / Reduce tiles number
  - Increase batch-size / Less patches / Low Overhead dependant)
  - Simple to implement / Reduce workflow complexity

- → Higher Accuracy
- → Lower Training Time (model
- → Improve Team productivity © Copyright IBM Corporation 2021

## **More informations**

- TensorFlow materiels
  - https://developer.ibm.com/linuxonpower/2021/06/11/tensorflow-large-model-supportresources/
  - TensorFlow Large Model Support Research Paper : <u>https://arxiv.org/pdf/1807.02037.pdf</u>
  - TensorFlow Large Model Support Case Study: <a href="https://developer.ibm.com/linuxonpower/2018/07/27/tensorflow-large-model-support-case-study-3d-image-segmentation/">https://developer.ibm.com/linuxonpower/2018/07/27/tensorflow-large-model-support-case-study-3d-image-segmentation/</a>
- IBM AC922 with NVLink 2.0 connections between CPU and GPU : <u>https://www.ibm.com/us-en/marketplace/power-systems-ac922</u>
- PyTorch Upstream contribution in progress
  - Code: <u>https://github.com/mtbrandy/pytorch/tree/v1.2.0-LMS</u>
  - Wiki: <u>https://github.com/mtbrandy/pytorch/wiki/Large-Model-Support</u>

### **Notices and disclaimers**

- © 2018 International Business Machines Corporation. No part of this document may be reproduced or transmitted in any form without written permission from IBM.
- U.S. Government Users Restricted Rights use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.
- Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information. This document is distributed "as is" without any warranty, either express or implied. In no event, shall IBM be liable for any damage arising from the use of this information, including but not limited to, loss of data, business interruption, loss of profit or loss of opportunity. IBM products and services are warranted per the terms and conditions of the agreements under which they are provided.
- IBM products are manufactured from new parts or new and used parts.

In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply."

 Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

- Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those
- customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.
- References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.
- Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.
- It is the customer's responsibility to insure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer's business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer follows any law.

### **Notices and disclaimers continued**

- Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products about this publication and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.
- The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

 IBM, the IBM logo, ibm.com and [names of other referenced IBM products and services used in the presentation] are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at: www.ibm.com/legal/copytrade.shtml.

•